

08/08/24

MAP

MLE

(Maximum likelihood estimate),

n data points,

$$D = \{x_1, x_2, \dots, x_n\}$$

likelihood of these data points, given a model  $\theta = \{ \}$ 

$$\log \hat{P}(D|\theta) = \hat{P}(\{x_1, x_2, \dots, x_n\}|\theta) = \prod_{k=1}^n \hat{P}(x_k|\theta)$$

the MLE of  $\theta$  is:  $\theta_{MLE} = \arg \max_{\theta} \hat{P}(D|\theta)$ Coin flip  
H | T

$$\theta = P(H) = q = \frac{\# \text{ of H}}{\text{total count}} = \frac{n_1}{n_1 + n_2}$$

$$D = \{x_1, x_2, \dots, x_n\}$$

$n_1: H$   
 $n_2: T$

$$P(D|\theta) = \prod_{i=1}^n P(x_i|\theta) = q^{n_1} (1-q)^{n_2}$$

to obtain  $\hat{q} = \arg \max_q q^{n_1} (1-q)^{n_2} = F$

$$\frac{\partial F}{\partial q} = n_1 q^{n_1-1} (1-q)^{n_2} - q^{n_1} (1-q)^{n_2-1} \times n_2 = 0$$

$$\Rightarrow q^{n_1-1} (1-q)^{n_2-1} (n_1(1-q) - q n_2) = 0$$

$$\Rightarrow q = \frac{n_1}{n_1 + n_2}$$

$$\log \hat{P}(D|\theta) = \sum_{i=1}^N \log(x_i|\theta) \quad D = \{x_1, \dots, x_n\}$$

MAP: Maximum a posteriori Estimate

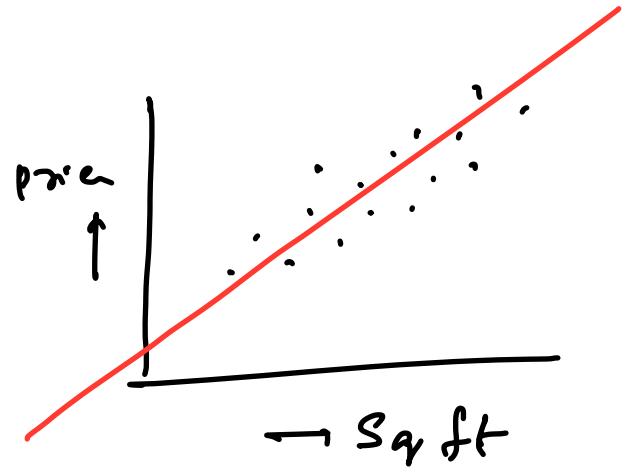
$$\log(\theta_{MAP}) = \arg \max_{\theta} P(\theta|D) = \arg \max_{\theta} \underbrace{P(D|\theta)}_{\text{likelihood}} \underbrace{P(\theta)}_{\text{prior}}$$

Posterior prob  $\nearrow$

$$= \log P(\theta) + \log P(D|\theta)$$
$$= \log P(\theta) + \sum_{i=1}^n \log P(x_i|\theta)$$

# Linear Regression

Living area (Sqft)	price (\$)
2104	450
1600	330
3000	232
→ 2500	?



$$h_0(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots$$

$$= \theta_0 + \sum_{i=1}^n \theta_i x_i$$

$$= \theta_0 + \sum_{i=0}^n \theta_i x_i$$

$x_0 = 1$

minimize the distance from all the points to the line  $L$

$$f(x_1, x_2, x_3) = x_1 + 3x_2 + 4x_3$$

(1, 2, 3)

$$f = \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 =$$

$$\text{Loss: } \frac{1}{2} \sum_{i=1}^n (f(x_i) - y_i)^2$$

$d_i = (x_i, y_i)$

$d_1$   
 $d_2$   
 $d_3$   
 $\vdots$   
 $d_m$

$$\theta = \{ \theta_0, \theta_1, \dots, \theta_n \}$$

$m = \#$  of training instance.

$$J(\theta) = \sum_{i=1}^m (h_0(x^i) - y^i)^2$$

Superscript  $\rightarrow$  instance

$j =$  Subscript  $\rightarrow$  parameter / feature

$$\frac{\partial J}{\partial \theta_j} = \frac{\partial}{\partial \theta_j} \frac{1}{2} (h_0(x^i) - y^i)^2$$

$$= 2 \times \frac{1}{2} (h_0(x^i) - y^i) \frac{\partial}{\partial \theta_j} (h_0(x) - y)$$

$$= (h_0(x^i) - y^i) \times \frac{\partial}{\partial \theta_j} \left( \sum_{k=1}^n \theta_k x_k - y \right)$$

$$= (h_0(x^i) - y^i) x_j$$

old  $\downarrow$

$\hat{\theta}_j \leftarrow \theta_j - \alpha \frac{\partial J}{\partial \theta_j}$

learning rate  $\alpha$

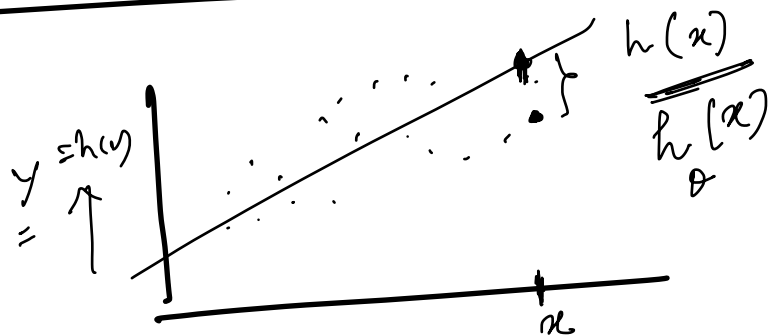
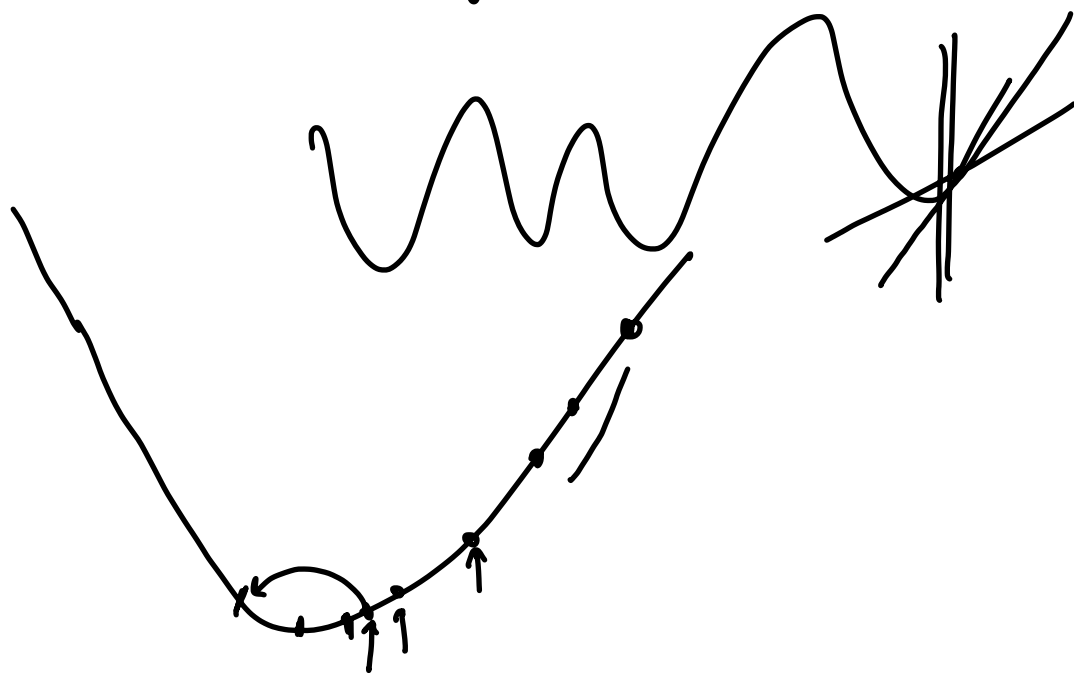
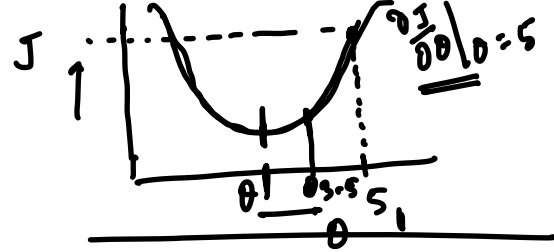
new  $\uparrow$

$$\frac{\partial J}{\partial \theta_p} = (h_0(x^i) - y^i) x_p$$

Batch GD | Stochastic GD | mini Batch GD

$$\hat{\theta}_j \leftarrow \theta_j - m \frac{\partial J}{\partial \theta_j}$$

$$\theta = 5$$



22/08/24

$$y = w_0 + w_1 x_1 + w_2 x_2 + \dots$$

$m = \#$  of data points  
 $n = \#$  of features

$$\frac{1}{2} \sum_{i=1}^m (h_0(x^i) - y^i)^2$$

$$x = \{ \dots \}_m$$

The Normal Equation

$$f: \mathbb{R}^{m \times m} \rightarrow \mathbb{R}$$

$$f(A) = \frac{3}{2} A_{11} + 5 A_{12}^2 + A_{21} A_{22}$$

$$A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix}$$

$$\nabla_A f(A) = \begin{bmatrix} \frac{\partial f}{\partial A_{11}} & \frac{\partial f}{\partial A_{12}} \\ \frac{\partial f}{\partial A_{21}} & \frac{\partial f}{\partial A_{22}} \end{bmatrix}$$

$$\text{tr}(A) = \text{tr} A = \sum_{i=1}^n A_{ii}$$

$$= \begin{bmatrix} 3/2 & 10 A_{12} \\ A_{22} & A_{21} \end{bmatrix}$$

①  $\text{tr}(ABC) = \text{tr}(CAB) = \text{tr}(BCA)$

②  $\text{tr}(A) = \text{tr}(A^T)$

④  $\text{tr}(aA) = a \text{tr} A$

③  $\text{tr}(A+B) = \text{tr}(A) + \text{tr}(B)$

⑤  $\nabla_A \text{tr}(AB) = B^T$       ⑥  $\nabla_{A^T} f(A) = (\nabla_A f(A))^T$

⑦  $\nabla_A \text{tr} ABA^T C = CAB + C^T A B^T$

$$\vec{X} = \begin{bmatrix} -x^{(1)T} \\ -x^{(2)T} \\ \vdots \\ -x^{(m)T} \end{bmatrix} \quad \vec{Y} = \begin{bmatrix} y^1 \\ y^2 \\ \vdots \\ y^m \end{bmatrix} \quad \vec{\theta} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}$$
  
 $m \times m$        $m \times 1$        $m \times 1$

$$\underline{(\vec{X}\vec{\theta} - \vec{Y})} = \begin{bmatrix} h_0(x^1) - y^1 \\ h_0(x^2) - y^2 \\ \vdots \\ h_0(x^m) - y^m \end{bmatrix}$$

$$J(\theta) = \frac{1}{2} \sum (h_0(x^i) - y^i)^2 = \frac{1}{2} (\underline{X\theta - Y})^T (\underline{X\theta - Y})$$

$$\nabla_{\theta} J(\theta) = \frac{1}{2} \nabla_{\theta} (\underline{X\theta - Y})^T (\underline{X\theta - Y})$$

$$= \frac{1}{2} \nabla_{\theta} (\theta^T X^T - Y^T) (X\theta - Y)$$

$$= \frac{1}{2} \nabla_{\theta} (\theta^T X^T X \theta - \theta^T X^T Y - Y^T X \theta + \underline{Y^T Y})$$

$$= \frac{1}{2} \left[ \nabla_{\theta} \text{tr} (\theta^T X^T X \theta - \theta^T X^T Y - Y^T X \theta) \right]$$

$$= \frac{1}{2} \left[ \nabla_{\theta} \text{tr} (\theta^T X^T X \theta) - \text{tr} (\theta^T X^T Y) - \text{tr} (Y^T X \theta) \right]$$

$$= \frac{1}{2} \left( \nabla_{\theta} \text{tr} (\theta^T X^T X \theta) - 2 \nabla_{\theta} \text{tr} (Y^T X \theta) \right)$$

$$= \frac{1}{2} \left( \nabla_{\theta} \text{tr} (\theta^T X^T X \theta) - 2 X^T Y \right)$$

$$\begin{aligned} \nabla_{\theta} \text{tr} (Y^T X \theta) &= \nabla_{\theta} \text{tr} (\theta Y^T X) \\ &= (Y^T X)^T \end{aligned}$$

$$\begin{aligned} A &= Y^T X \theta \\ A^T &= (Y^T X \theta)^T \\ &= \theta^T X^T Y \end{aligned}$$

$$\nabla_A \text{tr}(AB) = B^T$$

$$\nabla_{\theta} \text{tr} (\theta^T X^T X \theta)$$

$$\nabla_{\theta} \text{tr} \left( \begin{matrix} \theta & \theta^T & X^T & X \\ \hline A & B & A^T & C \\ \hline \vdots & \vdots & \vdots & \vdots \\ I & & & \end{matrix} \right)$$

$$= X^T X \theta + X^T X \theta = 2 X^T X \theta$$

$$\nabla_A \text{tr} ABA^T C = CAB + C^T A B^T$$

$$\downarrow \frac{1}{2} (2 x^T x \theta - 2 x^T y) = 0$$

$$\boxed{x^T x \theta = x^T y}$$

$$\theta = (x^T x)^{-1} x^T y$$

Normal Eq  
if  $x^T x$  is  
invertible

$$y = w_0 + w_1 x_1 + \dots + w_n x_n$$

$$= w_0 + w_1 x_1^2 + w_2 x_2^2 + \dots + w_n x_n^2 \leftarrow$$

linear in the parameters.

$$y = w_0 + w_1 \phi(x_1) + w_2 \phi(x_2) + \dots + w_n \phi(x_n)$$

$\phi$  = non-linear basis funt.

$$\phi_j(x) = x^j - \text{polynomial}$$

$$= \frac{(x - \mu)^2}{2\sigma^2} - \text{Gaussian}$$

$$= \frac{1}{1 + \exp(-s_1 x)} - \text{Sigmoid}$$

$$J = \sum_{i=1}^m (y^i - \sum_{j=1}^n \theta_j \phi(x^i))^2 = \sum_{i=1}^m (y^i - \theta^T \underbrace{\phi(x^i)}_{x^i})^2$$

$$\frac{\partial J}{\partial \theta} = \sum (y^i - \theta^T x^i)^2 = 2 \sum (y^i - \theta^T x^i) x^i = 0$$

$$\sum y^{(i)} x^{(i)T} = \theta^T \sum x^{(i)T} x^{(i)}$$

$$\Rightarrow x^T y = (x^T x) \theta \quad \left\| \begin{array}{l} \text{Normal} \end{array} \right.$$

$$\Rightarrow \boxed{\theta = (x^T x)^{-1} (x^T y)}$$

$$\underline{\underline{x x^T}}$$

$$(x^T x) = \text{non-invertible}$$

$$x \mid m \times n$$

Regularizer

Ridge

L2

Ridge Regmin

L1

Lasso Regmin

$$\hat{\theta} = \arg \min \sum_{i=1}^m (y_i - x_i \theta)^2 + \lambda \|\theta\|_2^2$$

$$\hat{\theta} = (X^T X + \lambda I)^{-1} X^T Y$$

$$\hat{\theta} = \sum_{i=1}^m (y_i - x_i \theta)^2 + \lambda \|\theta\|_1$$

Prove that:  $\nabla_A \text{tr}(ABA^T C) = C^T A B^T + C A B$

→ we use two theorems as follows:

$$\text{tr}(ABCD) = \text{tr}(DABC) = \text{tr}(CDAB) = \text{tr}(BCDA) \quad \text{--- (i)}$$

chain Rule of matrix differentiation:

$$d(AB) = d(A)B + A d(B) \quad \text{--- (ii)}$$

Therefore,  $\nabla_A \text{tr}(ABA^T C) = \nabla_X \text{tr}(X B A^T C) + \nabla_X \text{tr}(A B X^T C)$

$$\Rightarrow \nabla_X \text{tr}(X B A^T C) = (B A^T C)^T \quad (\because \nabla_A \text{tr}(AB) = B^T)$$

$$= C^T A B^T \quad \text{--- (iii)}$$

$$\nabla_X \text{tr}(A B X^T C) = \nabla_X \text{tr}(X^T C A B) \quad (\text{due to (i)})$$

$$\Rightarrow \left( \nabla_X \text{tr}(X^T C A B) \right)^T$$

$$\text{as } \nabla_{A^T} f(A) = \left( \nabla_A f(A) \right)^T \Rightarrow \left( \nabla_{X^T} \text{tr}(X^T C A B) \right)^T$$

$$\text{as } \nabla_A \text{tr}(AB) = B^T \Rightarrow \left( (C A B)^T \right)^T \Rightarrow C A B \quad \text{--- (iv)}$$

Then from (iii) and (iv)

$$\nabla_A \text{tr}(ABA^T C) = C^T A B^T + C A B$$

# Probabilistic Interpretation of Linear Reg.

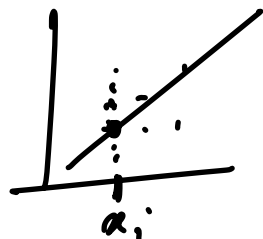
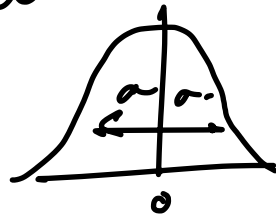
29/08/24

$$\underline{y^i} \neq \underline{\theta^T x^i} = h_{\theta}(x^i)$$

$$\underline{y^i} = \underline{\theta^T x^i} + \underline{\epsilon^i} \quad \text{error} \Rightarrow \text{unmodelled error}$$

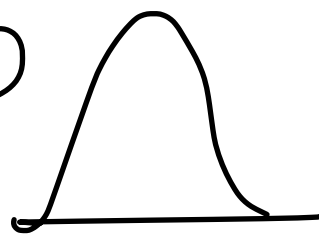
$$E(\epsilon^i) = 0$$

$$E(y^i) \approx E(\theta^T x^i)$$



$$p(\epsilon^i) = \frac{1}{\sqrt{2\pi} \sigma} \exp\left(-\frac{(\epsilon^i)^2}{2\sigma^2}\right)$$

$$\underline{p(y^i | x^i; \theta)} = \frac{1}{\sqrt{2\pi} \sigma} \exp\left(-\frac{(y^i - \theta^T x^i)^2}{2\sigma^2}\right) \quad \text{--- (1)}$$



$$y^i | x^i; \theta \sim \mathcal{N}(\theta^T x^i, \sigma^2)$$

$\epsilon \sim \text{iid}$

$$\underline{L(\theta)} = L(\theta; \vec{X}, \vec{Y}) = P(\vec{y} | \vec{X}; \theta) \quad \leftarrow$$

design matrix

vector of labels

$$= \prod_{i=1}^m p(y^i | x^i; \theta) = \prod_{i=1}^m \frac{1}{\sqrt{2\pi} \sigma} \exp\left(-\frac{(y^i - \theta^T x^i)^2}{2\sigma^2}\right) \quad \text{--- from (1)}$$

maximize

$$L(\theta) \equiv \max \log L(\theta) = \ell(\theta)$$

$$\ell(\theta) = \log L(\theta)$$

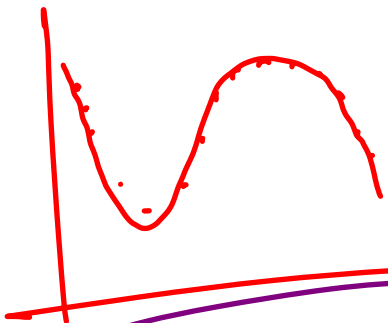
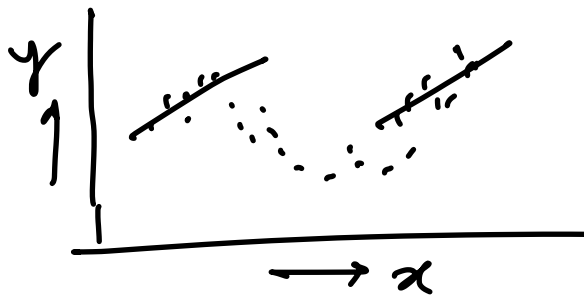
$$= \sum_{i=1}^m \log \left( \frac{1}{\sqrt{2\pi} \sigma} \exp\left(-\frac{(y^i - \theta^T x^i)^2}{2\sigma^2}\right) \right)$$

$$= m \log \frac{1}{\sqrt{2\pi} \sigma} - \frac{1}{\sigma^2} \sum_{i=1}^m \frac{(y^i - \theta^T x^i)^2}{2}$$

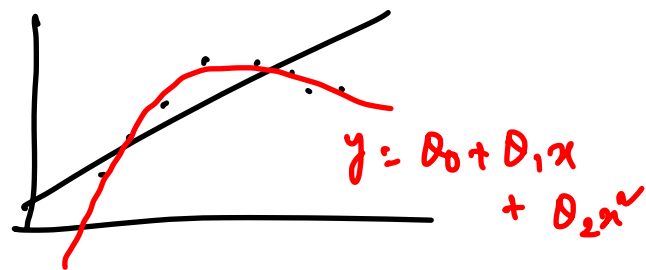
independent of  $\theta$

# Locally weighted LR

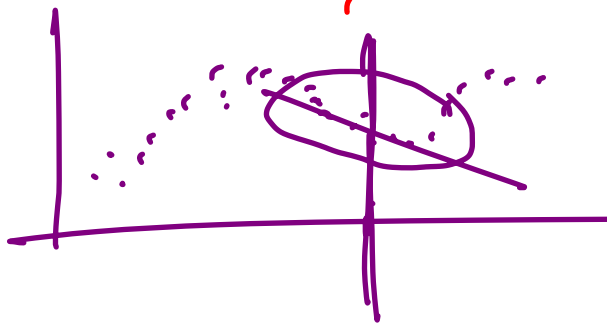
$$\begin{aligned}
 y &= \theta_0 + \theta_1 x \\
 &= \theta_0 + \theta_1 x + \theta_2 x^2 \\
 &= \theta_0 + \theta_1 x + \theta_2 x^2
 \end{aligned}$$



$$y = \sum_{i=1}^5 \theta_i x^i$$



Non-parametric model



Fit  $\theta$  by minimizing

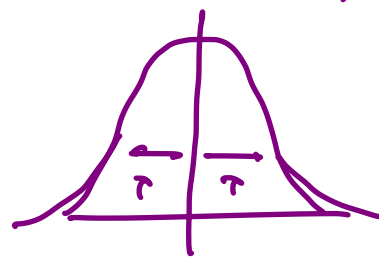
$$\sum_{i=1}^m w_i (y_i - \theta^T x_i)^2$$

$$w_i = \exp\left(-\frac{(x_i - x)^2}{2\tau^2}\right)$$

$\tau$  = bandwidth  
 = broad  
 = narrow

$\tau =$   
 over smooth thing / underfit

jagged fit / overfit



$SSE = SST - SSM \rightarrow$  Board

## Logistic Regression

31/08/24

- Classification method  $\rightarrow$  binary

$$y \rightarrow 0, 1 / -1, 1$$

0 / -1  $\rightarrow$  neg.

1 / +1  $\rightarrow$  positive.

$$x_i \rightarrow \{x_1, x_2, \dots, x_n\}$$

$$h_0(x) = \theta^T x = \theta_0 + \sum_{i=1}^n \theta_i x_i = \sum_{i=0}^n \theta_i x_i$$

$x_0 = 1$



$p(+)$        $p(-)$

$$h_{\theta}(x) = g(\theta^T x) = \frac{1}{1 + e^{-\theta^T x}} = \sigma(\theta^T x)$$

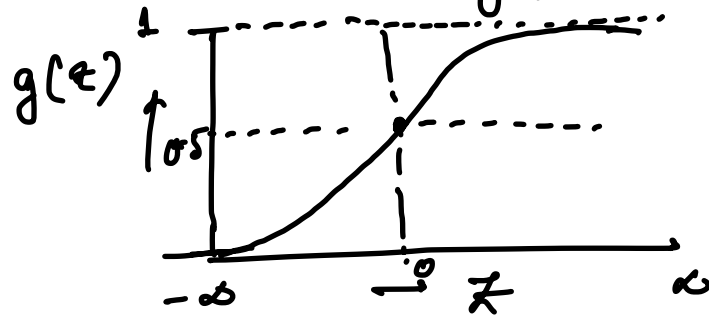
$$g(z) = \frac{1}{1 + e^{-z}}$$

→ logistic function  
Sigmoid fun<sup>n</sup>

$$g(z) \rightarrow 1; z \rightarrow \infty$$

$$g(z) \rightarrow 0; z \rightarrow -\infty$$

↓  
 $z=0$   
 $g = \frac{1}{2}$



$$\frac{\partial g}{\partial z} = \frac{\partial}{\partial z} \frac{1}{1 + e^{-z}}$$

$$= \frac{1}{(1 + e^{-z})^2} (e^{-z})$$

$$= \frac{1}{(1 + e^{-z})} \left( \frac{e^{-z}}{1 + e^{-z}} \right) = \frac{1}{1 + e^{-z}} \left( 1 - \frac{1}{1 + e^{-z}} \right)$$

$$= g(z) (1 - g(z))$$

$$g'(z) = g(z) (1 - g(z))$$

Assume that

$$p(y=1 | x^i; \theta) = h_{\theta}(x^i)$$

$$p(y=0 | x^i; \theta) = 1 - h_{\theta}(x^i)$$

$$x_i | i=1, \dots, m$$

$$p(\vec{y} | \vec{x}; \theta) = \prod_{i=1}^m p(y=1 | x^i)$$

$$\dots p(y=0 | x^i)$$

$$= \prod_{i=1}^m \left( h_{\theta}(x^i) \right)^{y^i} \left( 1 - h_{\theta}(x^i) \right)^{1 - y^i}$$

$$L(\theta) = p(\vec{y} | \vec{x}; \theta) = \prod_{i=1}^m h_{\theta}(x^i)^{y^i} (1 - h_{\theta}(x^i))^{1 - y^i}$$

maximize  $L(\theta)$

$$L(\theta) = \log L(\theta) = \sum_{i=1}^m y^i \log h_{\theta}(x^i) + (1-y^i) \log (1-h_{\theta}(x^i)) \quad (1)$$

$\theta \leftarrow \theta + \nabla_{\theta} L(\theta)$  - Gradient ascent

derivative of  $L(\theta)$  w.r.t.  $\theta = \{\theta_0, \theta_1, \theta_2, \dots, \theta_n\}$

$$\begin{aligned} \frac{\partial L(\theta)}{\partial \theta_j} &= y \frac{1}{g(\theta^T x)} g'(\theta^T x) \cdot \frac{\partial}{\partial \theta_j} (\theta^T x) - (1-y) \frac{1}{1-g(\theta^T x)} g'(\theta^T x) \frac{\partial}{\partial \theta_j} (\theta^T x) \\ &= \left( y \frac{1}{g(\cdot)} - (1-y) \frac{1}{1-g(\cdot)} \right) g'(\cdot) \frac{\partial}{\partial \theta_j} (\theta^T x) \\ &= \left( \frac{y}{g(\cdot)} - \frac{(1-y)}{1-g(\cdot)} \right) g(1-g) x_j \\ &= [y(1-g(\cdot)) - (1-y)g(\cdot)] x_j \\ &= (y - g(\cdot)) x_j = (y - h_{\theta}(x)) x_j \end{aligned}$$

$$\theta_j \leftarrow \theta_j + \alpha (y^i - h_{\theta}(x^i)) x_j^i \text{ for } x_i$$

why sigmoid

$$p(y=1|x) = \sum \theta_i x_i$$

odds  $\left| \frac{p(y=1|x)}{1-p(y=1|x)} \right| = \sum \theta_i x_i$

$$\Rightarrow \left( \frac{p}{1-p} \right) = e^{\sum \theta_i x_i}$$

$$\Rightarrow \frac{p}{1-p} = e^{\sum \theta_i x_i}$$

$$\Rightarrow p = (1-p) e^{\sum \theta_i x_i}$$

we proved on whiteboard that (1) is concave

$$\Rightarrow p = \frac{e^{\sum}}{1 + e^{\sum}}$$
$$p(y=1|x) = \frac{e^{\sum \theta_i x_i}}{1 + e^{\sum \theta_i x_i}}$$

$$p(y=0|x) = \frac{1}{1 + e^{\sum}}$$

---